

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representation of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

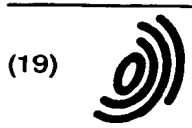
- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY

As rescanning documents *will not* correct images, please do not report the images to the Image Problem Mailbox.

This Page Blank (uspto)

09/557,282 #483



Europäisches Patentamt
European Patent Office
Office européen des brevets



(11) **EP 0 987 680 A1**

(12) **EUROPEAN PATENT APPLICATION**

(43) Date of publication:
22.03.2000 Bulletin 2000/12

(51) Int. Cl.⁷: **G10L 11/00**

(21) Application number: 99202980.1

(22) Date of filing: 13.09.1999

(84) Designated Contracting States:
**AT BE CH CY DE DK ES FI FR GB GR IE IT LI LU
MC NL PT SE**
Designated Extension States:
AL LT LV MK RO SI

(30) Priority: 17.09.1998 EP 98307574

(71) Applicant:
**BRITISH TELECOMMUNICATIONS public limited
company**
London EC1A 7AJ (GB)

(72) Inventor: **Marston, David Frank**
Bembridge, Isle Wight PO35 5SG (GB)

(74) Representative:
Nash, Roger William et al
BT Group Legal Services,
Intellectual Property Department,
Holborn Centre, 8th Floor,
120 Holborn
London EC1N 2TE (GB)

(54) **Audio signal processing**

(57) A speech coder (14) is operable to compress digital data representing speech using a Waveform Interpolation speech coding method. The coding method is carried out on the residual signal from a Linear Predictive Coding stage. On the basis of a series of overlapping frames of the residual signal, a series of respective spectra are found. The evolution of the spec-

tra is filtered in a multi-stage filtering process, the filtered phase data being replaced with the original phase data at the end of each stage. This is found to result in the decoder (28) being better able to approximate the original speech signal. This is of particular utility in relation to mobile telephony.

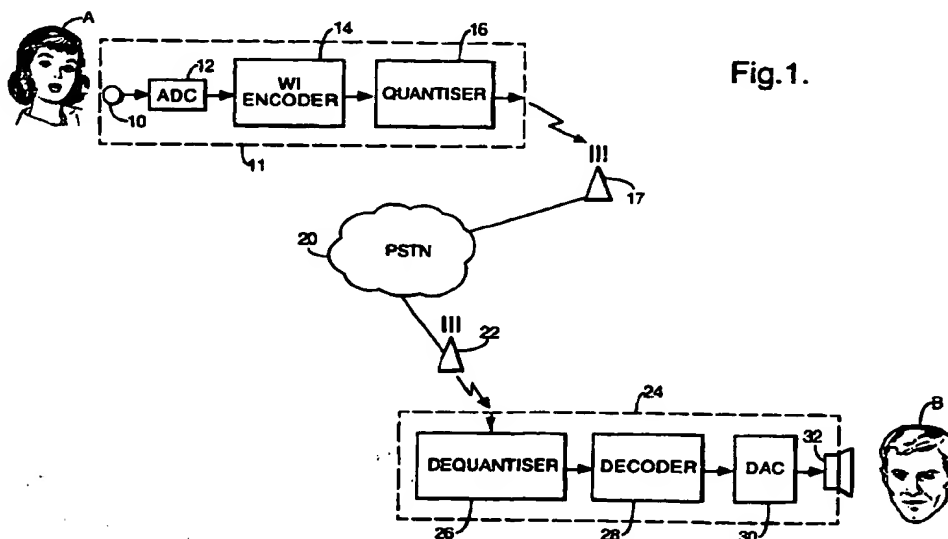


Fig.1.

EP 0 987 680 A1

Description

[0001] The present invention relates to audio signal processing. It has particular utility in relation to the separation of voiced speech and unvoiced speech in low bit-rate speech coders.

[0002] Low bit-rate speech coders are becoming increasingly commercially important as they enable a more efficient utilisation of the portion of the radio spectrum available to mobile phones.

[0003] Speech can be classified into three parts - voiced speech, unvoiced speech and silence. Any one of these may be corrupted by the addition of background noise. On a timescale of milliseconds, voiced speech can be viewed as a succession of repeated waveforms. This fact is exploited in a class of speech coding methods known as Prototype Waveform Interpolation (PWI) methods. Essentially, these methods involve sending information describing repeated pitch period waveforms only once, thereby reducing the amount of bits required to encode the speech signal. Initial PWI speech coding methods only encoded voiced speech, the other portions of the speech signal be coded using other methods (e.g. Code Excited Linear Prediction methods). One example of such a hybrid coding technique is described in "Encoding Speech Using Prototype Waveforms", W.B. Kleijn, IEEE Transactions on Speech and Audio Processing, Vol. 1, pp386-399, October 1993.

[0004] Later PWI methods were generalised so as to enable unvoiced speech and noise to be encoded as well. An example of such a method is described in "A General Waveform-Interpolation Structure for Speech Coding", W. B. Kleijn and J. Haagen, Signal Processing Theories and Applications, M. Hoit, C. Cowan, P. Grant, W. Sandham (Eds.), p1665-1668, 1994.

[0005] However, such coders have drawbacks in that the reconstituted speech sounds buzzy. The present inventors have established that the cause of this 'buzzy' is a poor separation of the voiced components of speech and the unvoiced/noisy components of speech.

[0006] According to a first aspect of the present invention there is provided a method of extracting one of a concordant component and a discordant component of a predetermined segment of an audio signal, said method comprising the steps of:

forming an initial evolution surface from a series of combined magnitude and phase spectra representing segments of said signal around said predetermined segment;

modifying said initial evolution surface to obtain a modified evolution surface representing said one of the concordant component or the discordant component of said signal; and

extracting said one of the concordant component or the discordant component of said predetermined segment from said modified evolution surface;

wherein said modifying step involves:

a plurality of component filtering steps and, prior to at least one of those filtering steps, the substitution of phase information derived from said initial evolution surface or an earlier one of the component steps for the phase information derived from the most recent component step.

[0007] Here, concordant is intended to refer to signals whose phase changes slowly in comparison to discordant signals whose phase changes more rapidly.

[0008] The present inventors have found that the rate of evolution of the phase information is useful in distinguishing between voiced speech (the concordant component of speech) and unvoiced speech/noise (the discordant component of speech).

[0009] However, it is likely that the invention will find application in other areas of audio signal processing such as the enhancement of noise-corrupted speech or music signals.

[0010] Conventional low-pass and high-pass Finite Impulse Response (FIR) digital filtering techniques do not reduce the magnitude of discordant and concordant signals respectively to zero. Therefore, they are limited in how well they can extract one of the concordant or discordant components of an audio signal.

[0011] A conventional FIR filter might be approximated by a series of shorter FIR filters. By decomposing a filtering process into a plurality of filtering stages and, in one or more of the intervals between those filtering stages, substituting phase information from an earlier stage for phase information from the most recent stage, a filtering process results which repeatedly uses the earlier phase information. Filtering a signal tends to smooth its phase and hence a filtered signal contains less information distinguishing its concordant and discordant parts. By reinstating the earlier phase information, the concordant or discordant component can be more thoroughly removed in the subsequent filtering stage(s). The result is a audio signal filtering process which is better able to extract a concordant or discordant component of an audio signal.

[0012] As suggested above, a repeated application of a low-pass filter will leave a modified evolution surface representing the concordant component of said predetermined segment. Preferably, each low-pass filtering step involves the application of an identical low-pass filter. This minimises the complexity of the processing method.

[0013] In preferred embodiments, the phase information derived from the initial evolution surface is used in all of said component steps. This maximises the effectiveness of the extraction method.

[0014] One way in which the discordant component can be calculated is to calculate the concordant component according to the first aspect of the present invention and subtract this from the original signal. Similarly, one way in which the concordant component can be calculated is to calculate the discordant component according to the first aspect of the present invention and subtract this from the original signal.

[0015] According to a second aspect of the present invention, there is provided an audio signal processor operable to extract one of a concordant component and a discordant component of a predetermined segment of an audio signal, said apparatus comprising:

means arranged in operation to form an initial evolution surface from a series of combined magnitude and phase spectra representing segments of said signal around said predetermined segment;

means arranged in operation to modify said initial evolution surface to obtain a modified evolution surface representing said one of the concordant component or the discordant component of said signal; and

means arranged in operation to extract said one of the concordant component or the discordant component of said predetermined segment from said modified evolution surface; wherein said apparatus further comprises:

means arranged in operation to carry out a plurality of filtering steps and, prior to at least one of those filtering steps, to substitute phase information derived from said initial evolution surface or an earlier one of the component steps for the phase information derived from the most recent component step.

[0016] According to a third aspect of the present invention, there is provided a speech coding apparatus including:

a storage medium having recorded therein processor readable code processable to encode input speech data, said code including:

initial evolution surface generation code processable to generate initial evolution surface data comprising combined magnitude and phase data for segments of said input speech data;

separation code processable to derive separate phase data and magnitude data from said input speech data;

evolution surface modification code processable to generate a modified evolution surface representing one of a voiced component or an unvoiced/noise component of said input speech data; and

component extraction code processable to extract said one of the voiced component or the unvoiced/noise component from said input speech data;

wherein said evolution surface modification code comprises:

evolution surface filtering code processable to filter said initial evolution surface data a plurality of times;

evolution surface decomposition code processable to derive magnitude data and phase data subsequent to one or more of said filtering steps; and

earlier phase reinstatement code processable to replace the phase data obtained on processing said evolution surface decomposition code with an earlier version of the phase data.

[0017] According to another aspect of the present invention there is provided a method of waveform interpolation speech coding comprising:

forming an initial evolution surface from a series of combined characteristic waveforms or spectra representing respective segments of said speech;

wherein said formation involves aligning each of said characteristic waveforms or spectra with an earlier characteristic waveform or spectrum of said series; and

said earlier waveform or spectrum is separated from the characteristic waveform or spectrum to be aligned with it by a variable number of members of said series, said variable number varying in accordance with the pitch of said signal.

[0018] It is found that the decoded version of unvoiced speech which has passed through a known waveform interpolation coder tends to have too high a periodic component. To reduce the undesirable periodic component in the output

version of unvoiced speech, alignment is made with a characteristic waveform or spectrum that is far enough back in the series to have a relatively low number of overlapping samples.

[0019] There now follows, by way of example only, a description of some embodiments of the present invention. The embodiments are described with reference to the accompanying drawings, in which:

Figure 1 is a schematic illustration of the application of a first embodiment of the present invention to a mobile telephony network;

Figure 2 shows the processes carried out in an encoder part of a mobile telephone forming part of the network of Figure 1;

Figure 3 is a schematic illustration of a spectral evolution surface produced during the operation of the encoder of Figure 2;

Figure 3B shows the evolution of an unvoiced speech frequency component over time;

Figure 3C shows the evolution of a voiced speech frequency component over time;

Figure 4 is a flow diagram which illustrates an evolution surface derivation method of prior-art encoders;

Figure 5 is a flow diagram which illustrates the evolution surface derivation method of the first embodiment;

Figure 6 shows the processes carried out by the decoder part of a mobile telephone according to the first embodiment of the present invention; and

Figure 7 illustrates the reduction of the unvoiced components of the evolution surface achieved using the method of the first embodiment.

[0020] A mobile telephone network (Figure 1) operating in accordance with a first embodiment of the present invention is operable to allow a first user A to converse with a second user B. User A's mobile phone is operable to transmit a radio signal representing parameters modelling user A's speech. The radio signal is received by a base station 17 which converts it to a digital electrical signal which it forwards to the Public Switched Telephone Network (PSTN) 20. The Public Switched Telephone Network 20 is operated to make a connection between base station 17 and a base station 22 currently serving user B. The digital electrical signal is passed across the connection, and, on receiving the signal, the base station 22 converts the digital electrical signal to parameters representing user A's speech. Thereafter, the base station 22 transmits a radio signal representing those parameters to user B's mobile phone 24. User B's mobile phone receives the radio signal and converts it back to an analogue electrical signal which is used to drive a loudspeaker 32 to reproduce A's voice. A similar communications path exists in the other direction from user B to user A.

[0021] For each of the radio communication sections, the mobile phone network selects an appropriate bit-rate for the parameters representing the user's speech from a full bit-rate (6.7kbits^{-1}), an intermediate bit-rate (4.6kbits^{-1}) and a half bit-rate (2.3kbits^{-1}).

[0022] The signal processing carried out in each mobile phone is now described in more detail. User A speaks into the microphone 10 of his mobile telephone 11 which converts his voice into an analogue electrical signal. This analogue signal is then passed to an Analogue to Digital Converter (ADC) 12 which digitises the signal to provide a 64kbits^{-1} digitally coded speech signal. A Waveform Interpolation (WI) encoder 14 receives the digitally coded speech signal and reduces it to a 6.7kbits^{-1} stream of parameters which represent user A's speech. The parameters are passed to a quantiser 16 which is operable to provide a variable rate parameter stream. The quantiser may simply forward the full-rate parameter stream or, if required, reduce the bit-rate of the parameter stream still further to the intermediate rate (4.6kbits^{-1}) or the half-rate (2.3kbits^{-1}).

[0023] It will be realised by those skilled in the art that the variable rate parameter stream undergoes further channel coding before being converted to a radio signal for transmission over the radio communication path to the base station 17.

[0024] User B's mobile phone recovers the variable rate parameter stream and, if required, uses interpolation to generate the 6.7kbits^{-1} parameter stream before passing the parameters to a decoder 28. The decoder 28 processes the parameter stream to provide a digitally coded reconstruction of user A's speech which is then converted to an analogue electrical signal by the Digital to Analogue Converter (DAC) 30, which signal is used to drive the loudspeaker 32.

[0025] The operation of the WI encoder 14 will now be described in more detail. The encoder 14 of user A's mobile phone receives the digitally coded speech signal from the Analogue to Digital Converter 30 and carries out a number

of processes (Figure 2) on the digitally coded speech signal to provide the stream of parameters representing user A's speech.

[0026] The encoder first divides the digitally coded speech signal into 10ms frames. Linear Predictive Coding (LPC) techniques (34,36,38) are then used in a conventional manner to provide, for each frame, a set of ten spectral shape parameters (Line Spectral Frequencies or LSFs) and a residual signal.

[0027] A pitch period detection process 40 provides a measure (expressed as a number of sample instants) of the pitch of the current frame of speech.

[0028] The residual signal then passes to a waveform extraction process which is carried out to obtain a characteristic waveform for each one of four 2.5ms sub-frames of each frame. Each characteristic waveform has a length equal to the pitch period of the signal at that sub-frame. Given that voiced speech normally has a pitch period in the range 2ms to 18.75ms, it will be realised that the characteristic waveforms will normally overlap one another to a significant degree. The residual signal for voiced speech has a sharp spike in each pitch period and the window used to isolate the pitch period concerned is movable by a few sample points so as to ensure the spike is not close to the edge of the window. Expressed in mathematical notation, the characteristic waveforms are obtained by windowing the residual signal as follows:

$$cw[i, k] = res(k + 20i - \frac{p_i}{2} - 1 - q) \text{ where } i = 0, 1, 2, 3 \text{ and } k = 1, 2, 3, \dots, p_i \quad \text{Equation 1}$$

[0029] Where $cw[i, k]$ represents the characteristic waveform for the i th sub-frame and $res(x)$ means the value of the x th sample of the residual signal. The pitch period from the pitch detector is p_i and, if required, q is increased from 0 to 4 in order to shift the spike in the residual away from the edge of the window.

[0030] The characteristic waveforms (of length p_i) thus extracted then undergo a Discrete Fourier Transform (DFT) 44 to produce, for each residual sub-frame, a characteristic spectrum. In mathematical notation, the characteristic spectra (CS) are calculated as follows:

$$CS[i, \omega] = \text{DFT}(cw[i, k], p_i) \quad \text{Equation 2}$$

[0031] Where $CS[i, \omega]$ is a complex value associated with a frequency interval ω and the i th sub-frame of the residual, the complex values for all frequency intervals forming a complex spectrum for the i th sub-frame of the residual. $cw[i, k]$ and p_i are as defined above.

[0032] The conventional technique of zero-padding is then used to expand the characteristic spectra so that they are all 76 values in length. To compensate for the effect of this on the power spectrum, the magnitude part of the characteristic spectra relating to shorter pitch periods is decreased in proportion to the pitch period associated with the residual sub-frame from which it is derived. In mathematical notation:

$$|CS_{norm}[i, \omega]| = \frac{150}{p_i} |CS[i, \omega]| \text{ where } \omega = 0, 1, 2, \dots, 76 \quad \text{Equation 3}$$

[0033] Where $|CS_{norm}[i, \omega]|$ represents the magnitude (or, in mathematical language, modulus) of the normalised complex spectral values and $|CS[i, \omega]|$ represents the magnitude of the complex value $CS[i, \omega]$ - p_i is as defined above.

[0034] It will be realised that the characteristic spectra are generally obtained from signal segments which overlap at least the signal segments used in deriving the previous and subsequent characteristic spectra. For voiced speech segments, there will be little difference in the magnitude of the complex values associated with each frequency interval of a spectrum and the corresponding magnitude values of the spectra derived from adjacent segments of the signal. However, the time offset between the adjacent signals manifests itself as a phase offset between adjacent spectra. In order to correct this phase offset the phase spectra (consisting of the phase, or, in mathematical language, argument of the complex spectral values) are operated on by alignment process 46.

[0035] Where the pitch period of the signal is long, a large number of samples may be used in calculating both a current spectrum and the spectra on either side. This leads to a similarity between adjacent spectra even in signals that are noisy in character. This similarity is undesirable since it reduces the distinction between voiced and unvoiced speech. In order to prevent such similarity arising in relation to unvoiced speech/noise each characteristic spectrum is aligned with another characteristic spectrum which may precede it by a many as four sub-frames. The interval (measured in sub-frames) between the characteristic spectra which are aligned with one another increases with increasing pitch period as follows:

if $p_4 < 90$ then $d = 1$
 if $90 \leq p_4 < 105$ then $d = 2$
 if $105 \leq p_4 < 125$ then $d = 3$
 if $p_4 \leq 125$ then $d = 4$

[0036] The alignment process shifts the phase values of one of the characteristic spectra to be aligned until the correlation between phase values of the two spectra reaches a maximum. The offset that is required to do this provides a phase correction for each one of the 76 frequency bins in the characteristic spectrum associated with a given sub-frame. The 'aligned' phase values are calculated by summing the original phase values and the phase correction (each is expressed in radians).

[0037] The phase spectrum is then combined with the magnitude spectrum associated with the sub-frame to provide an aligned characteristic spectrum for each sub-frame. Expressed mathematically,

$$CS_{aligned}[i, \omega] = |CS_{norm}[i, \omega]| e^{j \angle CS_{aligned}[i, \omega]} \quad \text{Equation 4}$$

[0038] Where j is $\sqrt{-1}$, and $\angle CS_{aligned}[i, \omega]$ represents the phase value obtained for the frequency interval ω associated with the i th sub-frame following the alignment procedure.

[0039] A normal representation of a spectrum has a series of bars spaced along a frequency axis and representing consecutive frequency intervals. The height of each bar is proportional to magnitude of the complex spectral value associated with the corresponding frequency interval. It is possible to visualise a further axis arranged perpendicularly to the frequency axis which represents the time at which a spectrum was obtained. Another spectrum derived a time interval later can then be visualised aligned with and parallel to the first spectrum and spaced therefrom in accordance with the scaling of the time axis. If this process is repeated for several spectra then a surface defined by the tops of the bars can be envisaged or computed from the individual magnitudes.

[0040] A simplified illustration of such a visualisation of the 'aligned' characteristic spectra output by alignment stage 46 is shown in Figure 3A (note that the alignment does not alter the magnitudes of the complex values forming the characteristic spectra and hence Figure 3A equally well represents the normalised characteristic spectra). For ease of illustration, only 11 spectral values are shown, rather than 76 as is actually the case in the embodiment.

[0041] The so-called 'evolution' of a spectral magnitude associated with a given frequency interval can be envisaged as the variation in that spectral magnitude over spectra derived from consecutive time intervals. The evolution of the magnitude associated with the second lowest frequency interval from time t_0 to t_4 in Figure 3A is therefore the succession of values V_1, V_2, V_3, V_4, V_5 .

[0042] As indicated above, the complex spectra in fact contain phase values as well as the magnitudes associated with a given frequency interval. The present inventors have found that an evolution of the complex spectral values associated unvoiced speech is more erratic than an analogous evolution derived from voiced speech. In particular, the phase component of the complex value varies more erratically for unvoiced speech. Figure 3B illustrates how a complex spectral value derived from unvoiced speech might evolve (the length of the line represents the magnitude, the angle α represents the phase). Figure 3C shows an evolution likely to be associated with voiced speech.

[0043] Returning to Figure 2, a Slowly Evolving Spectrum generation process 48 receives the aligned characteristic spectra and processes them to obtain a Slowly Evolving Spectrum. Conventionally, this has been done by storing, say, seven consecutive spectra and then applying a moving average filter to the evolution of the complex values associated with each frequency interval (Figure 4). Expressed in mathematical notation (here the complex spectral numbers are represented in the form of Real and Imaginary parts but the conversion from the Magnitude and Phase representation is trivial)

$$SES[i, \omega] = \sum_{m=-3}^3 a_m \operatorname{Re}\{CS_{aligned}[i+m, \omega]\} + j \sum_{m=-3}^3 a_m \operatorname{Im}\{CS_{aligned}[i+m, \omega]\} \quad \text{Equation 5}$$

[0044] Where $SES[i, \omega]$ represents the complex spectral values of a modified spectrum for the i th sub-frame of the residual signal and a_m represent the coefficients of the moving average filter.

[0045] According to the present embodiment, for each sub-frame, a series of operations are carried out on stored aligned characteristic spectra including the one associated with the current sub-frame and the six respectively associated with the six nearest sub-frames (Figure 5). In the first of these operations, a counter is set to zero (step 60). A moving average filter 62 is then applied to the evolutions of the complex spectral values associated with respective frequency intervals to provide a modified spectrum 64 to be associated with the current sub-frame.

[0046] The phase values of the modified spectrum are then replaced (step 66) by the phase values of the aligned characteristic spectrum associated with the current sub-frame to provide a hybrid characteristic spectrum 67 associated with the current sub-frame.

[0047] The counter is then increased by one (step 68) and a check is made on the value of the counter (step 70). If it has not yet reached six then the filtering 62 and phase replacement 66 steps are carried out on the hybrid characteristic spectrum just obtained.

[0048] If the counter has reached six then the magnitude values of the hybrid characteristic spectrum 67 obtained after the sixth replacement operation are output by the Slowly Evolving Spectrum generation process (Figure 2, 48) as the Slowly Evolving Spectrum 71 for the current sub-frame.

[0049] The Slowly Evolving Spectrum (SES) 71 is passed to the Rapidly Evolving Spectrum generation process 50. The Rapidly Evolving Spectrum (RES) generation process 50 subtracts the SES magnitude values from the corresponding magnitude values of the aligned characteristic spectrum associated with the current sub-frame to provide the magnitude values of the RES.

[0050] Both the SES magnitude values and the RES magnitude values are then arranged into Mel-scaled frequency intervals and the SES magnitude values 52 and RES magnitude values 54 for one out of every two sub-frames are forwarded to the quantiser (Figure 1, 16).

[0051] As explained in relation to Figure 1, the stream of parameters (pitch 41, RES magnitude values 54, SES magnitude values 52, LSFs 37) output by the WI encoder 14 are received at the decoder 28 in user B's mobile phone 24.

[0052] The processes carried out in the decoder 28 are now described with reference to Figure 6. The SES magnitude values 52 are passed to a phase generation process 80 which generates phase values to be associated with the magnitude values on the basis of known assumptions. In this embodiment the phase values are generated in the way described in the Applicant's International Patent Application No. PCT/GB97/02037 published as WO 98/05029. The phase values and the SES magnitude values are combined to provide a complex SES characteristic spectrum.

[0053] The RES magnitude values 54 are combined with random phase values 82 to generate a complex RES characteristic spectrum.

[0054] Interpolation processes 84,86 are then carried out on the two types of spectra to obtain one spectra of each type every 2.5ms. The two spectra thus created are then combined 88 to provide an approximation to a characteristic spectrum for each sub-frame. The approximate characteristic spectrum is then passed, together with the pitch 41, to a cubic interpolation synthesis process 90 which operates in a known manner to reconstruct an approximation to the residual signal originally derived in the LSF analysis process in the encoder (Figure 2, 38). A filter 92 which is the inverse of the analysis filter (Figure 2, 38) is then used to provide an approximation of the audio signal originally passed to the encoder (Figure 1, 14).

[0055] Owing to the nature of the decoding process (Figure 6), it is important that the SES magnitude values 52 are very low for unvoiced speech. If this is not the case then unvoiced speech components are synthesised in the same way as voiced components which results in the output speech sounding buzzy.

[0056] In the above-described embodiment of the present invention, the SES generation process (Figure 2, 48) is better able to reduce the SES magnitude values associated with unvoiced speech than the processes used in prior-art PWI encoders. In prior-art coders the erratic evolution of the phase values does result in the low-pass filtering operation (Figure 4, 57) reducing the magnitude values of the resultant SEW for the corresponding frequency interval. However, the present invention improves on this since it gives extra weight to the phase information in the characteristic spectrum (it will be recalled that it is the phase information that especially distinguishes unvoiced speech/noise from voiced speech). Extra weighting of the phase information is achieved by replacing the phase values at each stage of the iterative filtering process and thereby reintroducing the erratic phase values that particularly distinguish voiced and unvoiced speech before the next filtering stage. The result is low SES magnitudes associated with unvoiced speech and hence a less buzzy output than known encoders.

[0057] The reduction of the SES magnitude with repeated filtering stages is illustrated in Figure 7. It can be seen that there is little reduction in magnitude values associated with voiced speech, but that repeated iterations of the filter strongly reduce the magnitude values associated with unvoiced speech.

[0058] In other embodiments, in the SES generation process, the phase values obtained from any earlier filtering stage could be used to replace the phase resulting after a later filtering stage. Such a method would still provide a degree of improvement over the prior-art.

[0059] The above described processes (40, 42, 44, 46) which extract SES magnitude values from the residual signal could be used to derive a voicing measure for each of the frequency bands for each sub-frame. The voicing measure might simply be the ratio of the output SES magnitude to the original characteristic spectrum magnitude for a given frequency interval. Such a set of processes might be useful in a Multi-Band Excitation speech coder.

[0060] At the expense of extra processing, the alignment stage 46 might be included within the repeated processes contained within the loop illustrated in Figure 5. This would correct any drift introduced by the filtering process.

[0061] Those skilled in the art will be able to conceive of many different low-pass filters that may be used in the low-

pass filtering process 62.

[0062] In the above embodiment, each of the characteristic spectra corresponds to a single pitch period of the residual signal. Instead, the characteristic waveforms could be of a fixed length allowing the use of an efficient Fast Fourier Transform (FFT) algorithm to calculate the characteristic spectra. The characteristic spectra might then contain peaks and troughs corresponding to the fundamental of the input signal (which, of course, need not be a residual signal). The application of the iterative process described in relation to Figure 5 would then retain the peaks but reduce the troughs further. Such a method is likely to have application in noise reduction algorithms that might be applied to speech, music or any other at least partly periodic audio signals.

[0063] The improved separation of the spectra representing the unvoiced and voiced speech might also find application in speech recognition devices.

Claims

1. A method of extracting one of a concordant component and a discordant component of a predetermined segment of an audio signal, said method comprising the steps of:
 - forming an initial evolution surface from a series of combined magnitude and phase spectra representing segments of said signal around said predetermined segment;
 - modifying said initial evolution surface to obtain a modified evolution surface representing said one of the concordant component or the discordant component of said signal; and
 - extracting said one of the concordant component or the discordant component of said predetermined segment from said modified evolution surface;

wherein said modifying step involves:

a plurality of component filtering steps and, prior to at least one of those filtering steps, the substitution of phase information derived from said initial evolution surface or an earlier one of the component steps for the phase information derived from the most recent component step.
2. A method according to claim 1 wherein said component steps comprise respective low-pass filtering steps whereby said modification step provides a modified evolution surface representing the concordant component of said predetermined segment.
3. A method according to claim 2 wherein each low-pass filtering step involves the application of an identical low-pass filter.
4. A method according to any preceding claim wherein phase information derived from said initial evolution surface is used in all of said component steps.
5. A method according to any preceding claim further comprising the step of calculating the other of the concordant component and the discordant component by subtracting said one of the two components from said initial evolution surface.
6. A method according to claim 1 wherein said component steps comprise respective high-pass filtering steps whereby said modification step provides a modified evolution surface representing the discordant component of said predetermined segment.
7. A method according to claim 1 wherein said audio signal is substantially periodic and each predetermined segment represents a different pitch period.
8. A method of separating voiced speech from unvoiced speech and noise, said method comprising the steps of any preceding claim where said audio signal represents speech and said voiced speech corresponds to said concordant component and said unvoiced speech and noise corresponds to said discordant component.
9. A method of speech coding comprising the separation method of claim 8 whereby more information is used to code the voiced speech than is used to code the unvoiced speech and noise.

10. An audio signal processor operable to extract one of a concordant component and a discordant component of a predetermined segment of an audio signal, said apparatus comprising:

5 | means arranged in operation to form an initial evolution surface from a series of combined magnitude and phase spectra representing segments of said signal around said predetermined segment;

 | means arranged in operation to modify said initial evolution surface to obtain a modified evolution surface representing said one of the concordant component or the discordant component of said signal; and

10 | means arranged in operation to extract said one of the concordant component or the discordant component of said predetermined segment from said modified evolution surface;

 | wherein said apparatus further comprises:

15 | means arranged in operation to carry out a plurality of filtering steps and, prior to at least one of those filtering steps, to substitute phase information derived from said initial evolution surface or an earlier one of the component steps for the phase information derived from the most recent component step.

11. A speech coding apparatus including:

20 | a storage medium having recorded therein processor readable code processable to encode input speech data, said code including:

 | initial evolution surface generation code processable to generate initial evolution surface data comprising combined magnitude and phase data for segments of said input speech data;

25 | separation code processable to derive separate phase data and magnitude data from said input speech data;

 | evolution surface modification code processable to generate a modified evolution surface representing one of a voiced component or an unvoiced/noise component of said input speech data; and

 | component extraction code processable to extract said one of the voiced component or the unvoiced/noise component from said input speech data; wherein said evolution surface modification code comprises:

30 | evolution surface filtering code processable to filter said initial evolution surface data a plurality of times;

 | evolution surface decomposition code processable to derive magnitude data and phase data subsequent to one or more of said filtering steps; and

 | earlier phase reinstatement code processable to replace the phase data obtained on processing said evolution surface decomposition code with an earlier version of the phase data.

- 35 12. A method of waveform interpolation speech coding comprising:

 | forming an initial evolution surface from a series of combined characteristic waveforms or spectra representing respective segments of said speech;

40 | wherein said formation involves aligning each of said characteristic waveforms or spectra with an earlier characteristic waveform or spectrum of said series; and

 | said earlier waveform or spectrum is separated from the characteristic waveform or spectrum to be aligned with it by a variable number of members of said series, said variable number varying in accordance with the pitch of said signal.

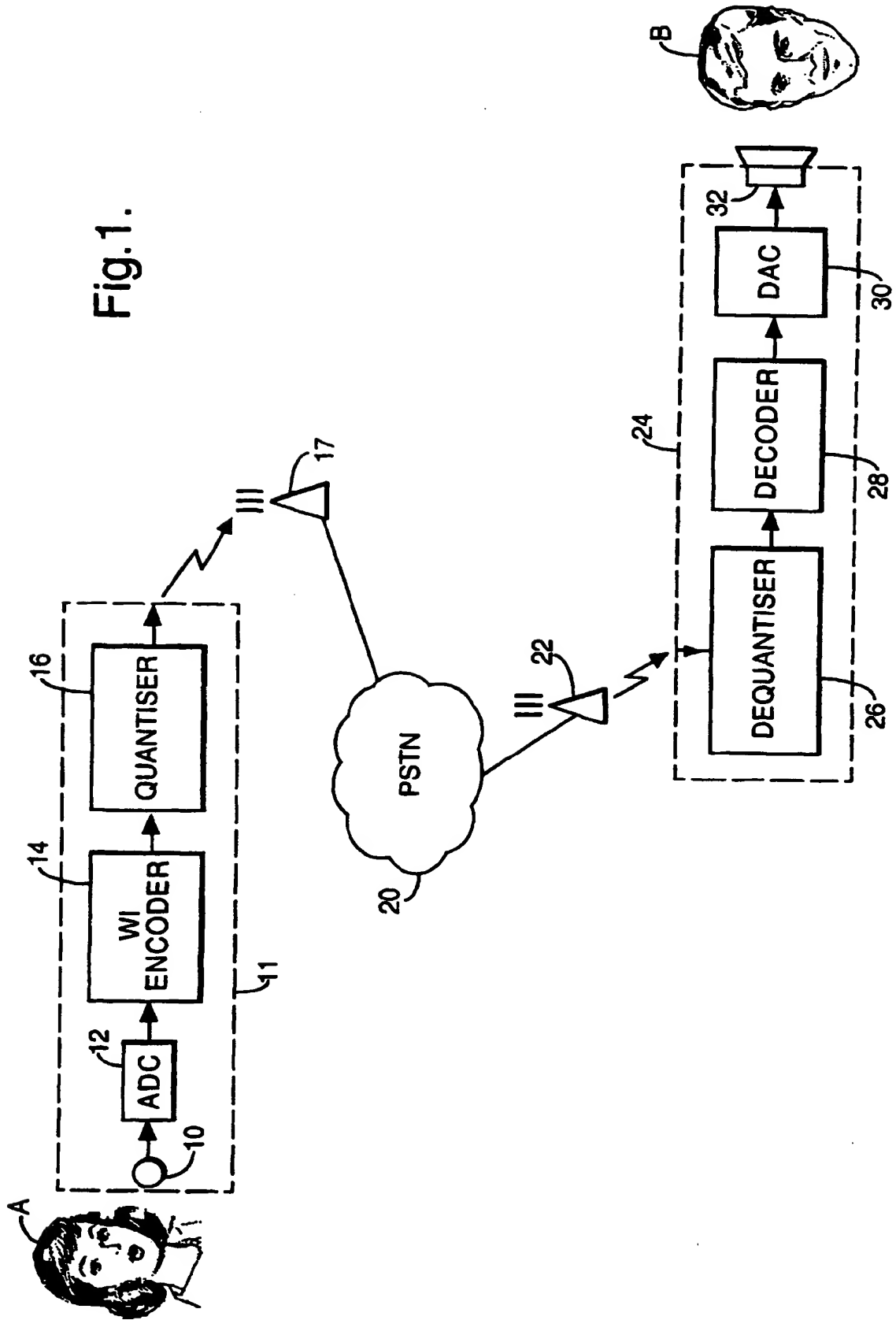
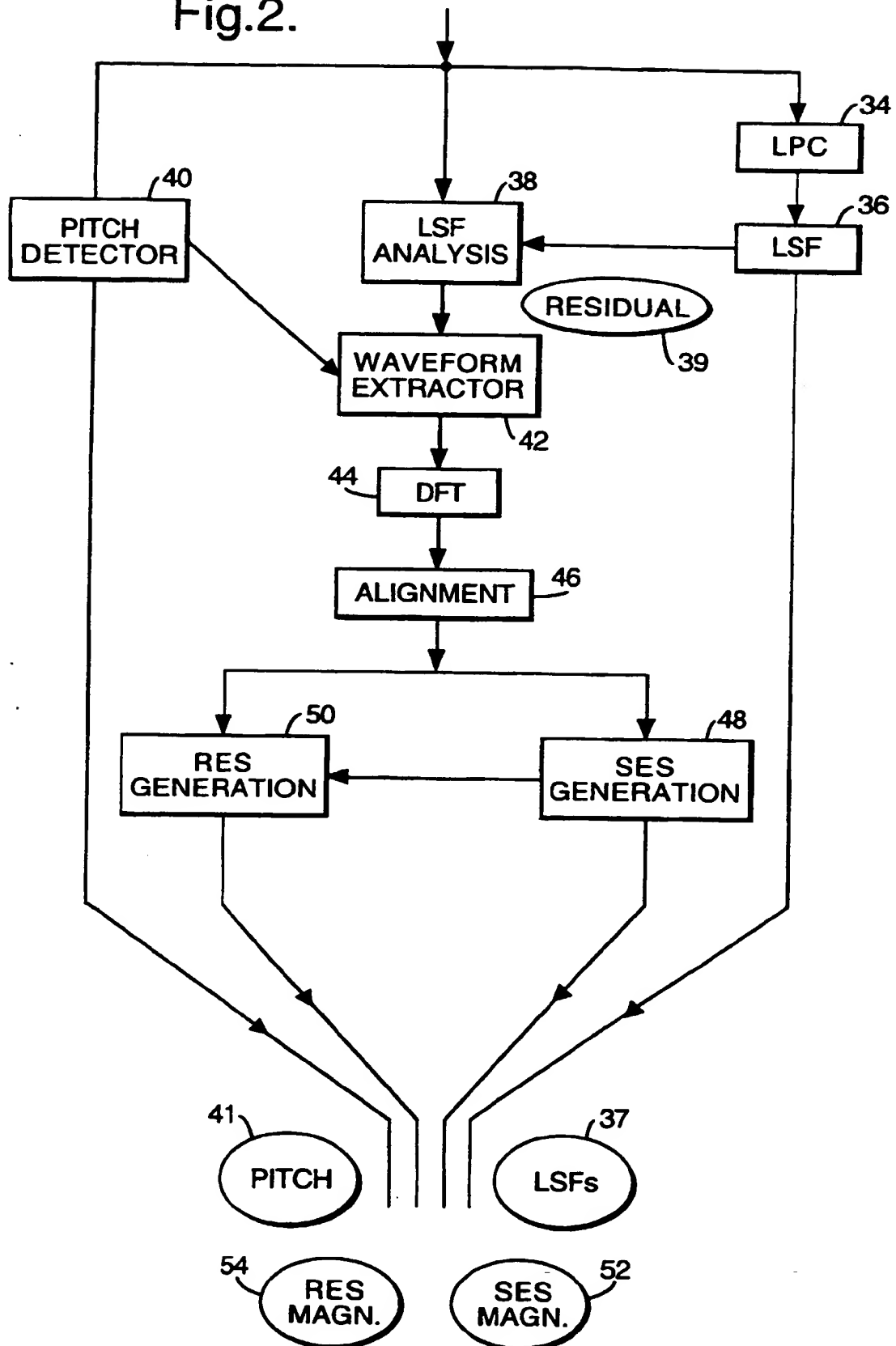


Fig. 1.

Fig.2.



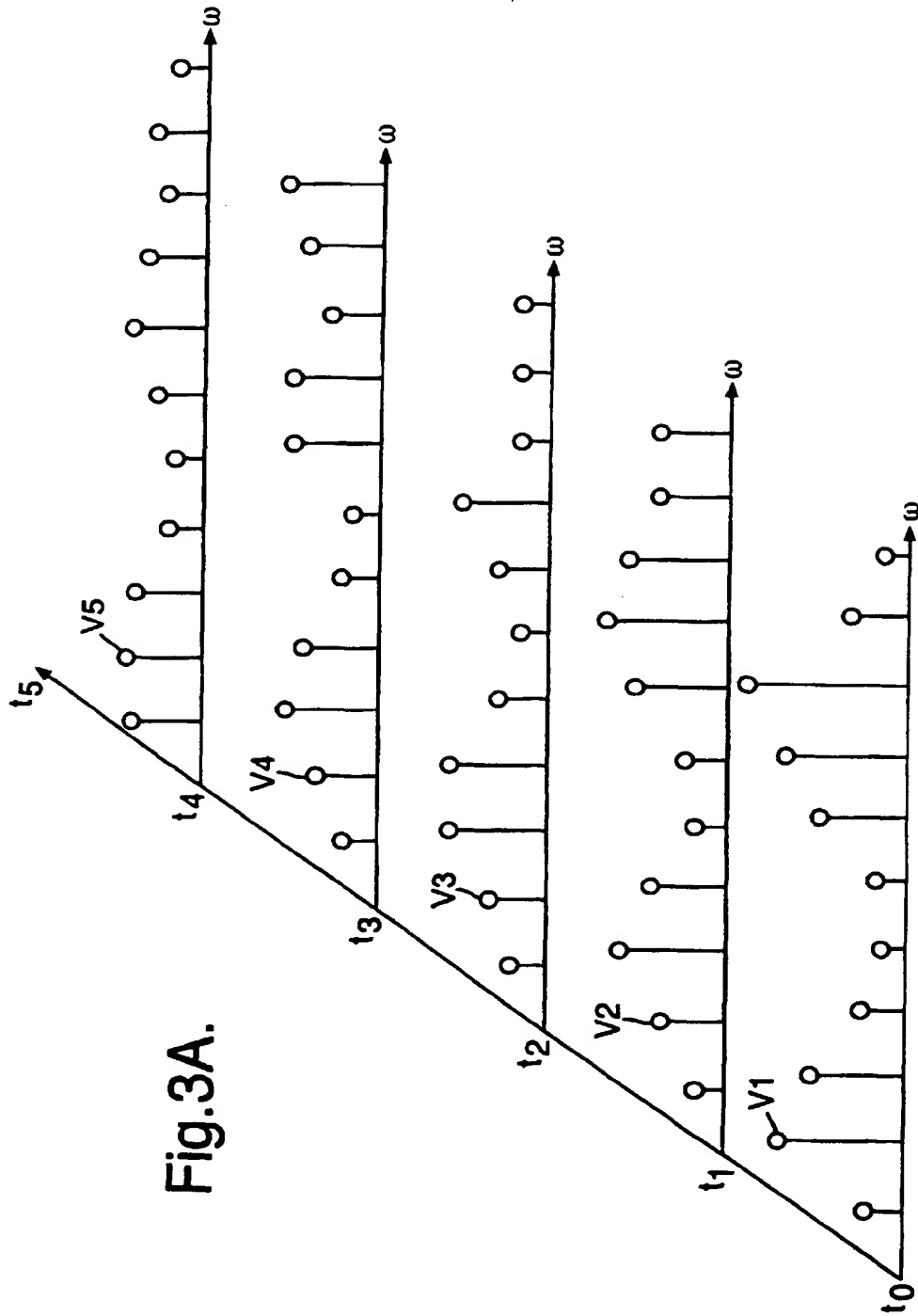


Fig. 3A.

Fig.3B.

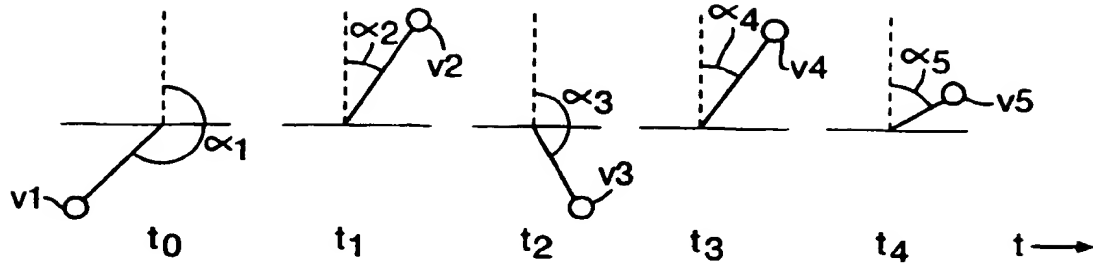


Fig.3C.

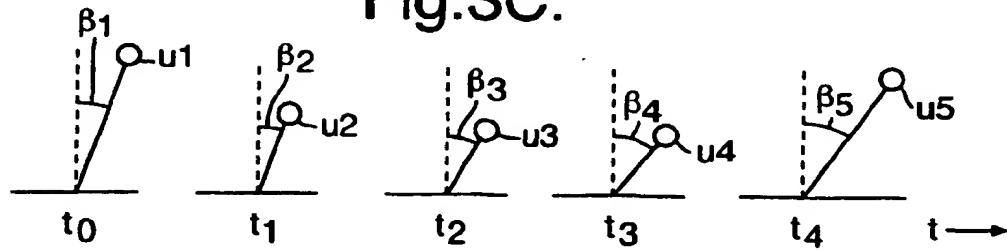


Fig.4.

PRIOR ART

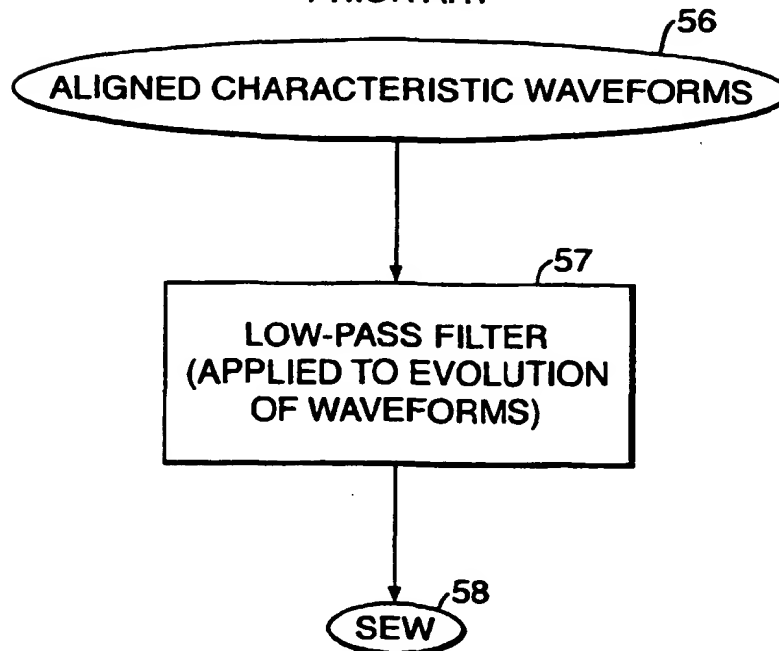
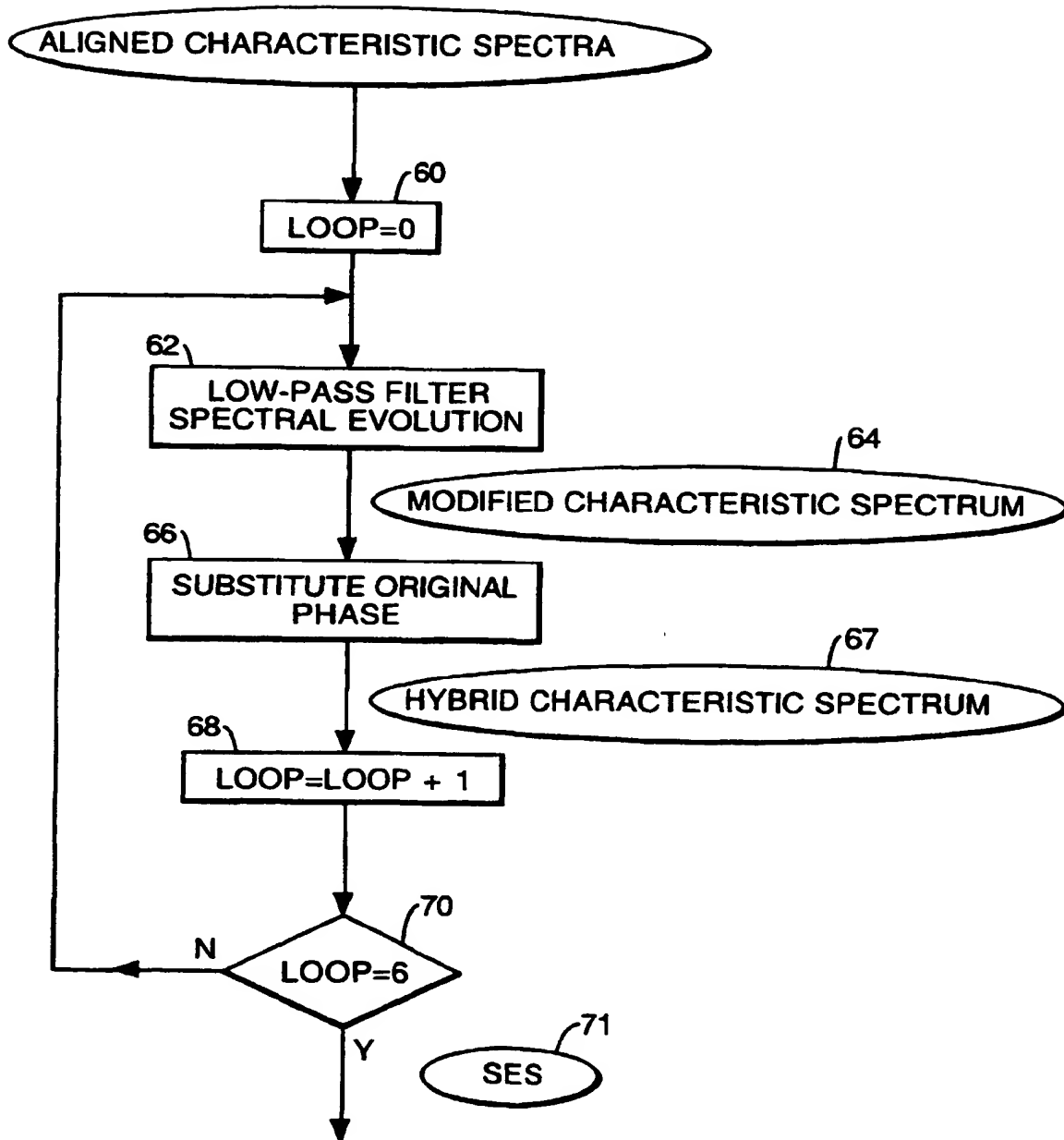


Fig.5.



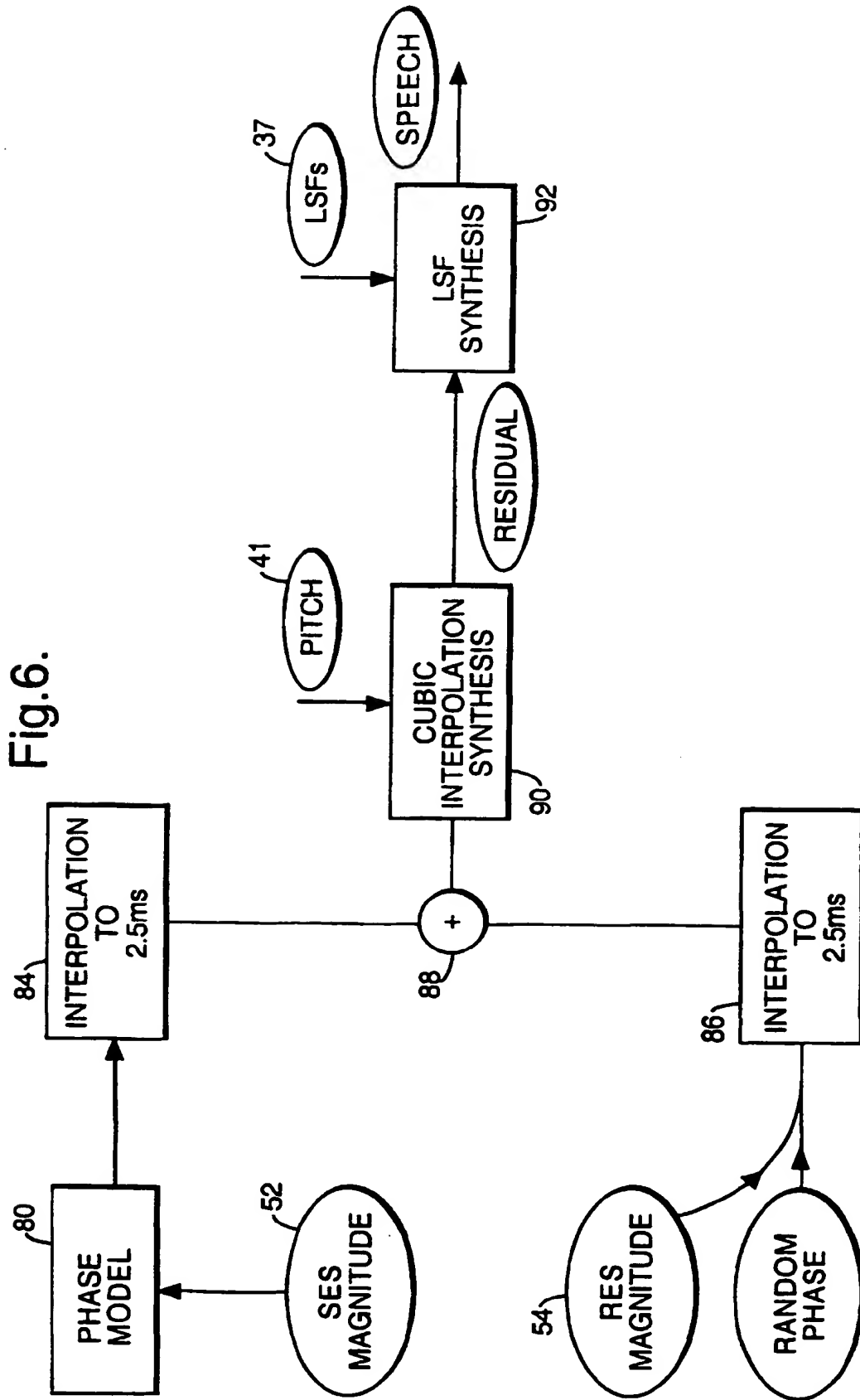
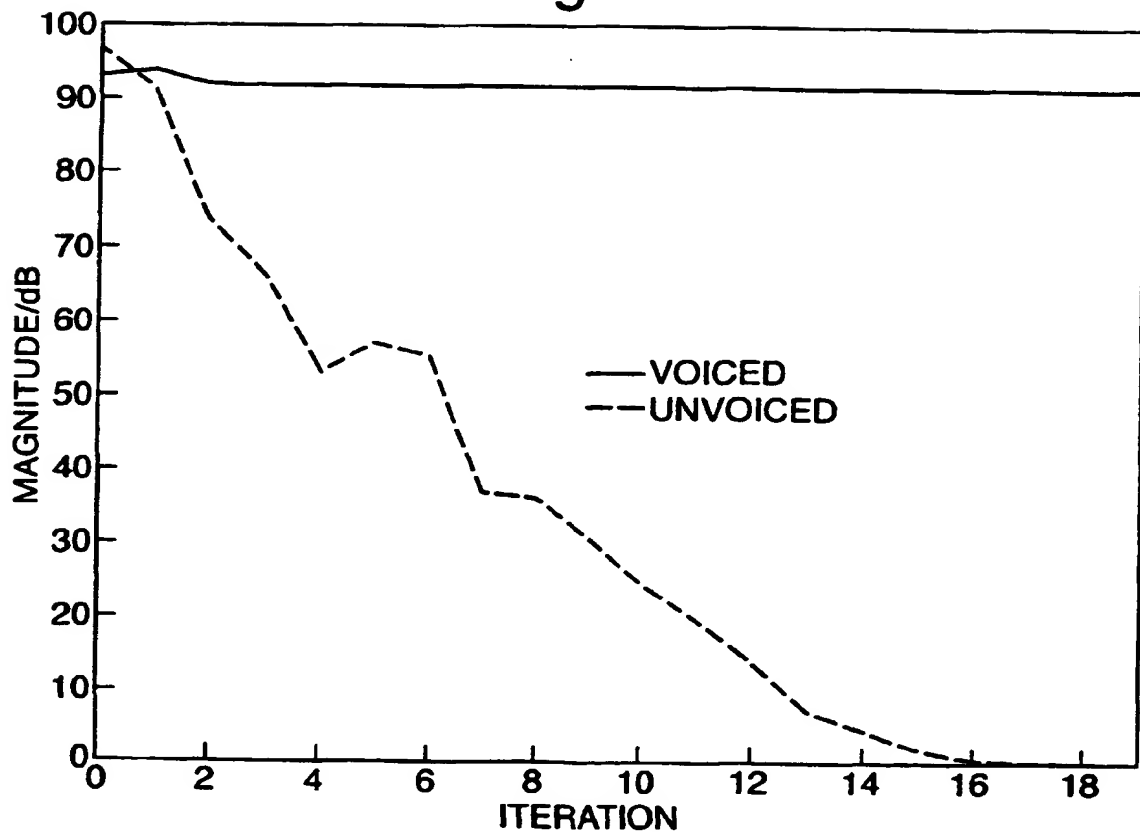


Fig.7.





European Patent
Office

EUROPEAN SEARCH REPORT

Application Number
EP 99 20 2980

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (InCL7)
A	EP 0 666 557 A (AT & T) 9 August 1995 (1995-08-09) * page 2, line 53 - page 3, line 21 * * page 10, line 15 - line 25 *	1,10-12	G10L3/00
A	EP 0 865 029 A (LUCENT) 16 September 1998 (1998-09-16) * page 3, line 21 - line 36 *	1,10-12	
A	WO 98 05029 A (CHEETHAM ET AL.) 5 February 1998 (1998-02-05) * abstract *	1,10-12	
A	CHONG ET AL.: "Use of the pitch synchronous wavelet transform as a new decomposition method for WI" PROCEEDINGS OF THE 1998 IEEE INTERNATIONAL CONFERENCE ON ACOUSTICS, SPEECH AND SIGNAL PROCESSING, ICASSP '98, vol. 1, 12 - 15 May 1998, pages 513-516, XP002093343 SEATTLE, WA, US	1,10-12	
D,A	KLEIJN: "ENCODING SPEECH USING PROTOTYPE WAVEFORMS" IEEE TRANSACTIONS ON SPEECH AND AUDIO PROCESSING, vol. 1, no. 4, 1 October 1993 (1993-10-01), pages 386-399, XP000422852 New York, NY, US * paragraph '000I! - paragraph '00II! * -/-	1,10-12	TECHNICAL FIELDS SEARCHED (InCL7) G10L
The present search report has been drawn up for all claims			
Place of search THE HAGUE		Date of completion of the search 20 October 1999	Examiner Lange, J
CATEGORY OF CITED DOCUMENTS X: particularly relevant if taken alone Y: particularly relevant if combined with another document of the same category A: technological background O: non-written disclosure P: intermediate document T: theory or principle underlying the invention E: earlier patent document, but published on, or after the filing date D: document cited in the application L: document cited for other reasons &: member of the same patent family, corresponding document			

EPO FORM 1503 (04/98) (P04001)



European Patent
Office

EUROPEAN SEARCH REPORT

Application Number
EP 99 20 2980

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (IntCL7)
D,A	<p>KLEIJN ET AL.: "A speech coder based on decomposition of characteristic waveforms" PROCEEDINGS OF THE INTERNATIONAL CONFERENCE ON ACOUSTICS, SPEECH, AND SIGNAL PROCESSING (ICASSP 1995), vol. 1, 9 - 12 May 1995, pages 508-511, XP000658042 DETROIT, MI, US</p> <p>* paragraph '0001! - paragraph '0003! *</p>	1,10-12	
			<p>TECHNICAL FIELDS SEARCHED (IntCL7)</p>
<p>The present search report has been drawn up for all claims</p>			
Place of search		Date of completion of the search	Examiner
THE HAGUE		20 October 1999	Lange, J
<p>CATEGORY OF CITED DOCUMENTS</p> <p>X : particularly relevant if taken alone Y : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document</p> <p>T : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application L : document cited for other reasons</p> <p>& : member of the same patent family, corresponding document</p>			

EPO FORM 1503 03.92 (P0401)

**ANNEX TO THE EUROPEAN SEARCH REPORT
ON EUROPEAN PATENT APPLICATION NO.**

EP 99 20 2980

This annex lists the patent family members relating to the patent documents cited in the above-mentioned European search report. The members are as contained in the European Patent Office EDP file on
The European Patent Office is in no way liable for these particulars which are merely given for the purpose of information.

20-10-1999

Patent document cited in search report	Publication date	Patent family member(s)	Publication date
EP 0666557 A	09-08-1995	US 5517595 A	14-05-1996
		CA 2140329 A	09-08-1995
		JP 7234697 A	05-09-1995
EP 0865029 A	16-09-1998	US 5924061 A	13-07-1999
		DE 69800011 D	02-09-1999
		JP 10319996 A	04-12-1998
WO 9805029 A	05-02-1998	AU 3702497 A	20-02-1998
		EP 0917709 A	26-05-1999

EPO FORM P449

For more details about this annex : see Official Journal of the European Patent Office, No. 12/82

This Page Blank (uspto)